



crh

L'ÉCOLE  
DES HAUTES  
ÉTUDES EN  
SCIENCES  
SOCIALES

# Les entités nommées (et au-delà) comme outil d'analyse de corpus textuels : une perspective de projet en humanités numériques

Carmen Brando (PGHN / CRH, EHESS)  
carmen.brand@ehess.fr

# Les entités nommées (EN) en traitement automatique des langues (TAL)

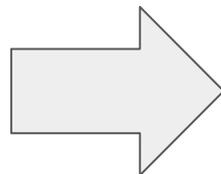
Extrait de David Copperfield de Charles Dickens

...  
I was born at Blunderstone, in  
Suffolk, or 'there by', as they say  
in Scotland.

...  
I remarked that,  
once or twice when Mr. Quinion  
was talking, he looked at Mr.  
Murdstone sideways

...

texte brut



...  
I was born at **Blunderstone**, in  
**Suffolk**, or 'there by', as they  
say  
in **Scotland**.

...  
I remarked that,  
once or twice when **Mr. Quinion**  
was talking, he looked at **Mr.**  
**Murdstone** sideways

...

Texte annoté en EN

# Les EN en TAL

- ❖ “*Il s’agit de types d’unités lexicales particulières qui font référence à une entité du monde concret dans certains domaines spécifiques notamment humains, sociaux, politiques, économiques, et qui est désignée par son nom*” (Ester)
- ❖ associée aux **noms propres** (*Barack Obama, Hugo, rue Mouffetard*) et parfois à des **descriptions définies** (un groupe nominal, ex : *le chat noir*)
- ❖ Les trois **catégories** (dites *coarse-grained*) : noms de personnes, de lieux, d’organisations ; extension de la typologie d’EN pour inclure : les fonctions de personnes, les productions humaines, les expressions numériques, temporelles, etc.

# Les EN en TAL: types d'ambigüités

- ❖ le même nom est utilisé pour plusieurs entités
  - *Paris* (France) et *Paris* (Texas)
- ❖ une même entité peut avoir plusieurs noms
  - *Paris*, *Paname*
- ❖ le nom d'une entité peut être utilisé pour désigner une entité en catégories différentes (le cas de la métonymie)
  - la *Sorbonne*, la *France*, ... une organisation mais peut être aussi un lieu selon le contexte
  - “Le *prix Nobel de la Paix* s'est montré digne devant une telle épreuve”

# Structurer et enrichir des corpus numériques/numérisés : avec les EN et au delà

## Annuaires des propriétaires et de propriétés de Paris (PTM)

### LISTE ALPHABÉTIQUE DES RUES DE PARIS

N.B. Les rues sont classées suivant l'ordre officiel de la Ville de Paris.

(Voir l'appendice à la fin du volume).

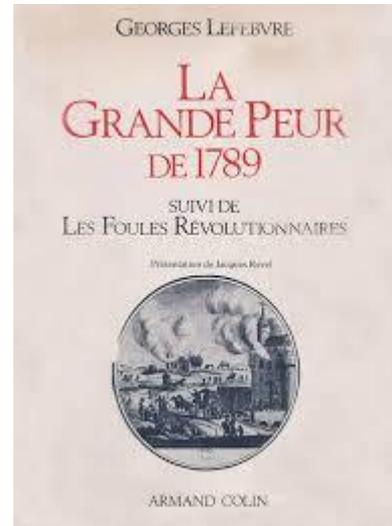
- A.9 ABBAYE (Rue de 1)<sup>14</sup>**
- 1 Entrée pass. Petite-Boucherie, 1.
  - 3 à 7 Lagulouërie (C. de), Paris, pl. St-Germain-des-Prés, 3.
  - 11-13 Vallès, Paris, r. Abbaye, 13.
  - 17-19 Entrée r. Rennes, 44.
  - 2 Entrée r. Rohaudé, 18.
  - 20-4 Entrée r. Furstenberg, 7-9.
  - 6 Michas et Douaette, Paris, r. Denfert-Rochereau, 47.
  - 8 Delmas, Paris, r. Lata, 2.
  - 10 Peignot, Paris, boul. Eug. Quinet, 68.
  - 12 Romazy, Vincennes, av. Marigny, 37.
  - 14 St-Luc (C. de), Paris, r. Abbaye, 14.
  - 16 Langre (Vve), Poissy, av. Mignat, 16 (S. et O.).
  - 18 Entrée r. Bonaparte, 37.
  - 20 Entrée r. Bonaparte, 42.
  - 22 Boudrouil (de), Paris, r. Bonaparte, 29.
  - 24 Entrée r. St-Jacques, 11.
- A.3 ABBÉ DE L'ÉPÉE (Rue de 1)<sup>11</sup>**
- 1 Mothe (M<sup>me</sup>), Paris, r. Abbé de l'Épée, 1.
  - 3 Madot, Sicy, y Sedan (H<sup>is</sup>-Sp.).
  - 5 Mahille, Paris, r. Abbé de l'Épée, 5.
  - 7 Saullier, Paris, r. Abbé de l'Épée, 7.
  - 9 Bricault, Paris, r. St-Antoine, 222.
  - 11 Entrée r. St-Jacques, 254.
  - 2 Entrée r. Guy-Lassas, 48.
  - 4 Yaver, Paris, r. Abbé de l'Épée, 4.
  - 6 Gérard, Thionny (S.-et-M.).
  - 8 Paris, Paris, r. Abbé de l'Épée, 8.
  - 10 Ville de Paris.
  - 12 Mossier (M<sup>me</sup>), Paris, r. Abbé de l'Épée, 12.
  - 14 Wickmann, gér. Delail, r. Abbé de l'Épée, 14.
  - 16 Fabre (H<sup>is</sup>), gér. Delail, Paris, r. Abbé de l'Épée, 16.
  - 18 Lortel (Vve), Paris, r. Mont-Thabor, 6.
  - 20 Entrée boul. St-Michel, 105.
- A.4 ABBÉ GRÉGOIRE (Rue de 1)<sup>11</sup>**
- 1 Gurdson, Paris, r. Abbé Grégoire, 1.
  - 3 Breton (Vve), Paris, r. Abbé Grégoire, 3.
  - 3<sup>bis</sup>-5 Breton, Paris, boul. St-Germain, 182.
- 7 Guillaumes, Paris, r. Henard, 29.
  - 9 St-Euze (M<sup>me</sup>), gér. St-Marie, Paris, r. Mademo, 33.
  - 11 Marens, Paris, fg St-Honoré, 164.
  - 13 Pamart, Paris, r. Abbé Grégoire, 13.
  - 15 Bot (Vve), Paris, r. Abbé Grégoire, 15.
  - 17 Entrée r. Clerche-Midi, 64.
  - 19 Entrée r. Clerche-Midi, 57.
  - 21 Hamot, Douaette (S.-et-O.).
  - 23 Costot (Vve), Chalon, av. Pâtes d'Bas, 12 (S.-et-O.).
  - 25 Herivieux (Vve), Paris, r. Abbé Grégoire, 25.
  - 27 Reulat, Paris, r. Le Peletier, 42.
  - 29 Lemaux (H<sup>is</sup>), gér. Roche-Faye, Paris, r. Geruelle, 90.
  - 31 Meunier, Paris, r. Rivoli, 5.
  - 33 Lefèvre, Paris, r. Abbé Grégoire, 33.
  - 35 Dubois, Paris, r. Abbé Grégoire, 35.
  - 37 Roux, Paris, boul. Montparnasse, 89.
  - 39 Falmouier (H<sup>is</sup>), gér. M<sup>me</sup> Dapuy, Vincennes, r. Paris, 56 (Seine).
  - 41 Jansou, Paris, r. Abbé Grégoire, 41.
  - 43 Rollier (M<sup>me</sup>), St-Mandé, av. Hocheillon, 43 (Seine).
  - 45 Dubouché, Paris, r. Montmartre, 15.
  - 47 Entrée r. Vaugrain, 90.
  - 2 Minot, Coulommiers, r. Pvs, 13 (S.-et-M.).
  - 4 Bricault, Paris, boul. Montparnasse, 51.
  - 6 Dally, Paris, r. Lénie, 3.
  - 8 à 10 Dames de St-Maur, Paris, r. Abbé Grégoire, 10.
  - 18 Entrée r. Clerche-Midi, 66.
  - 20 Talbot, St-Maurice (Seine).
  - 22 David, Paris, r. Caumartin, 37.
  - 24 Quicrux, Paris, r. Abbé Grégoire, 24.
  - 26 Targat, Paris, av. République, 19.
  - 28 Entrée r. Vaugrain, 92.
- A.5 ABBÉ GROULT (Rue de 1)<sup>11,12</sup>**
- 1 Tréjier, Paris, r. Abbé Groault, 1.
  - 3 Entrée r. Rottembourg, 108.
  - 5<sup>5<sup>me</sup></sup> Avianin (H<sup>is</sup>), gér. Hallain, Paris, r. Ecole, 36.
  - 7 Fallot, Paris, r. Abbé Groault, 7.
  - 9 Jaeger, Paris, r. Puvie, 17.

## ::Savoirs

LE FIL DES IDÉES

Parcourez les textes de Savoirs

Titre	Auteur(s)
Alphabète et procédés de chat un savoir genévain du XVII <sup>e</sup> siècle	Jean-François Bait
Activités perceptives et cognitives : le cadastre sarde (XVIII <sup>e</sup> siècle)	Françoise Bregeot
Culture japonaise et culture numérique : du Web au Gutenberg des médias du XXI <sup>e</sup> siècle	Marc Jobling
Éthique du design biomimétique	Christophe Jacob
Les conclusions de savoirs ou temps des Bouquins combinateurs (Vie-les-lettres ou...)	Jean-Louis
Gestes et formes de l'écriture savante	Boisot Mandressi
La table de lettres dans la Chine impériale	Richard Schneider
L'écriture des mathématiciens	Marin Andler
Un siècle européen : la correspondance de Franz Curtius	Caroline Brocard
L'univers et le web : l'exposition internationale japonaise de 2005	Sophie Soubart
Réseaux, géolocalisation, contrats, collectifs, savoirs	Boisot Mandressi
Structures et éprouvés du savoir à distance : le cas de l'exploration robotique en ligne	Emmanuel Benveniste, Nil
La mobilité entre universités au Moyen Âge	Ad Tancop
Un théâtre de papier	Édith Fliediger
Les « paroles » du lettré : la construction humaniste d'un instrument de lecture	Jean-Marc Chatalekin
Lettrés, savoirs, savoirs	Nicolas Andler
Les échelles du savoir	Jean-Jacques Glissant



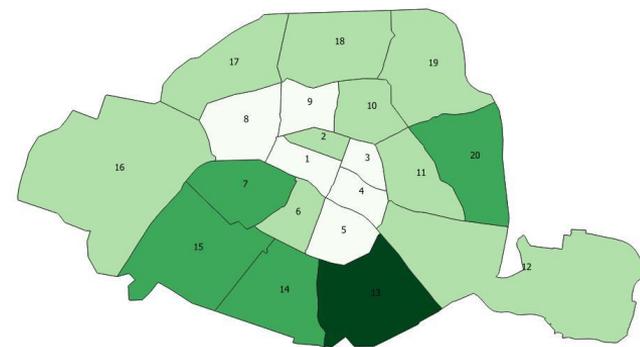
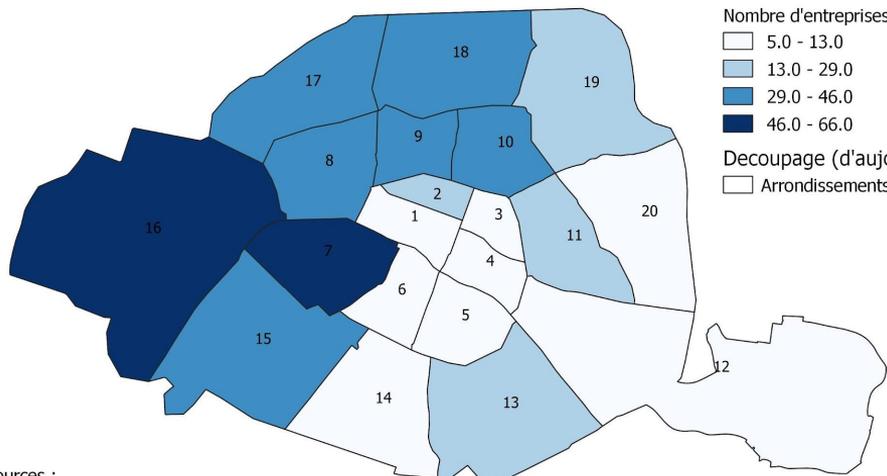
## Dictionnaire géo-historique d'Alcedo sur l'Amérique hispanique (ANR TopUrbi)

## La grande peur de 1789 de Georges Lefebvre

# Applications des EN dans les HN

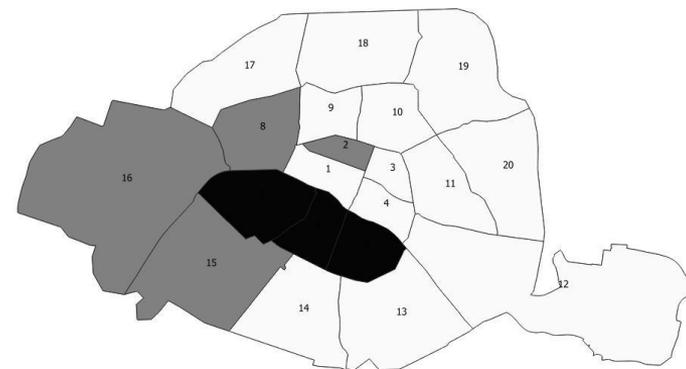
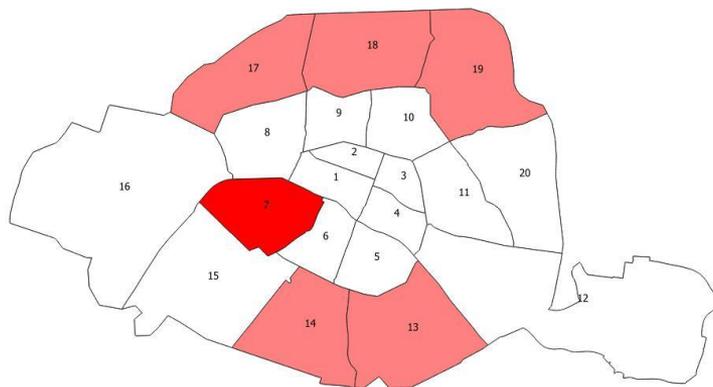
- Cartographier et visualiser les EN d'un corpus, par exemple, une cartographie ou analyse de la dimension spatiale et spatio temporelle d'un corpus
- Interroger de manière croisée des corpus et bases de données différents reliés par des EN
- Indexer pour préparer une édition numérique enrichie
- Calculer de suggestions de lecture de textes dans un corpus

## La propriété des entreprises en 1898 à Paris



### Sources :

- adresses : Annuaire de propriétaires de Paris et de la Seine, vol. 1898. Les localisations sont issues de <https://geo.api.gouv.fr/adresse>.
- découpage administratif : Ville de Paris (<https://opendata.paris.fr/>)



# Mise en oeuvre de chaînes de traitement des sources aux analyses

- tout en intégrant des approches issues des **différents domaines** : TAL, analyse des données, géomatique
- plusieurs approches/outils et solutions partielles existent mais l'enjeu est de faciliter (**généraliser**) leur **réutilisation** dans des projets avec des besoins proches
- Adapter la pipeline selon les **niveaux de structuration des textes** ; les EN, une composante de la pipeline
- Produire des données **réutilisables et exploitables** pour des recherches, les enrichir et spatialiser
- développer et **systematiser** des analyses à partir de ces données
- Pipeline EN/TAL vs pipeline HN : objectifs proches mais différents (unitâche vs. multitâche)

# Mise en oeuvre d'une chaîne de traitement des sources aux analyses : sur les outils

- performance des outils
- robustesse et généricité des outils
- outils en production avec une communauté et un appui technique
- interfaces utilisateurs accessibles à des utilisateurs non experts
- formats standards d'échange divers et variés

- 68-70 Huraud, Paris, av. Ternes, 38.  
 72 Devauchelle et Longuet, Bouchon Somme.  
 74 Clouzier, Paris, r. Montmartre, 157.  
 76 Gillet, Paris, r. N.-D. Nazareth, 36.  
 78 Renault, Noisy-le-Grand (S.-et-O.).  
 80 Boroux, Paris, boul. Beaujour, 43.  
 82 Grimbert, Paris, cité Talma, 11.  
 84 Daguet (M<sup>re</sup>), Paris, r. Mahillon, 18.  
 86 Bertron-Auger, La Flèche (Sarthe).  
 88 Calmels, Paris, boul. St-Marcel, 9.  
 90 Deville, Paris, r. Beuxelles, 23.  
 92 Decolange, Versailles, av. St-Cloud, 20.  
 94 Cros, Paris, r. Lombards, 23.  
 96 Masson (M<sup>re</sup>), gér. Rouchon, Paris, r. Quincampoix, 96.  
 98 Vincent, Paris, r. Quincampoix, 98.  
 100 Sanoner (Vve), Paris, r. Rivoli, 130.

#### 0.9 QUINZE-VINGTS (Passage des) <sup>45</sup>

- 1 Entrée r. Lyon, 46.  
 3 Davenne, Paris, pass. Quinze-Vingts, 3.  
 s. n° Ville de Paris.

#### R.1 RABELAIS (Rue) <sup>30,31</sup>

- 1 Riffel, Paris, r. Rabelais, 1.  
 3 Delouze, Paris, r. Rabelais 3.  
 5 Entrée r. Montaigne, 26.  
 2-4 Gérard (B<sup>re</sup>), Paris, fg St-Honoré, 85.  
 6-8 Aligre (M<sup>re</sup> d'), Paris, fg St-Honoré, 89.

#### R.4 RACINE (Rue) <sup>22</sup>

- 1 Ferry, Paris, r. Turin, 27.  
 3 Levoux (M<sup>re</sup>), Paris, r. Rennes, 108.  
 5 Clous, gér. Ferrembach, Paris, r. Racine, 5.  
 7 Bourrières, gér. Arnould, Paris, r. Racine, 7.  
 9 Barois, St-Germain-en-Laye (S.-et-O.).  
 11 Ville de Paris.  
 13 Finelle, St-Mandé, Ch<sup>te</sup> Etang, 12 (Seine).  
 15 Gille (H<sup>re</sup>), Paris, r. Racine, 15.  
 17 Chanroux, gér. Roussel, Paris, r. Racine, 17.  
 19 Entrée r. Monsieur le Prince, 22.  
 21 Normand (Vve), Paris, boul. Beaumarchais, 50.  
 23-25 Blondeau (M<sup>re</sup>), Paris, av. Percier, 8<sup>bis</sup>.  
 2 Crochard (M<sup>re</sup>), gér. Vve David, Paris, r. Médicis, 7.  
 4-6 Ville de Paris.  
 8 Ecole Nationale des Arts décoratifs (État).  
 s. n° Ecole de Médecine.  
 24 Entrée r. Monsieur le Prince, 20  
 26-28 Flammarion, St-Mandé, av. Daumesnil, 39 (Seine).  
 30 Croville, Paris, r. Monsieur le Prince, 2

#### R.6 RADZIWILL (Rue) <sup>3</sup>

- 9 Raincourt (M<sup>re</sup>), St-Germain-en-Laye, r. Voltaire, 44 (S.-et-O.).  
 11 Blot (Vve), gér. Amelin, Paris, r. St-Petersbourg, 24.  
 13 Croix, Paris, r. Puteaux, 7.  
 15-17 Fattet (M<sup>re</sup>), Paris, av. Champs-Élysées, 21.  
 19 Bertet, Paris, r. Radziwill, 19.  
 21 Charles (Vve), Paris, gal. Montpensier, 25.  
 23 Delâtre, Paris, r. Washington, 30.  
 25-27 Banque de France.  
 29 Chevet (Vve), Paris, r. François Gérard, 43.  
 31 Chauvelot, gér. Viron et Provost, Paris, r. Radziwill, 31.  
 33-35 Banque de France.  
 37 Cousin, Paris, r. Boccador, 3.  
 2 Banque de France.

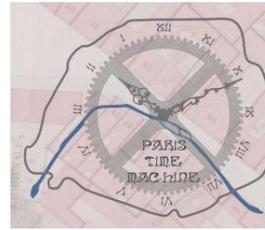
#### R.7 RAFFET (Rue) <sup>61</sup>

- 1 Entrée r. Source, 36.  
 3-5 Thomas (Vve), Paris, r. Beaujour, 49.  
 s. n° Ville de Paris.  
 19 R'chly, Paris, r. Raffet, 19.  
 19<sup>bis</sup> Soloville, Paris, boul. Beaujour, 4.  
 21 Hallot (M<sup>re</sup>), Paris, r. Raffet, 21.  
 23 St-Lanne, Paris, r. Raffet, 23.  
 s. n° Robert, Paris, pl. Poissou, 2.  
 s. n° Noblet, Paris, r. Cujas, 13.  
 2 Prévost, St-Germain-en-Laye, r. Paris, 2 (S.-et-O.).  
 4-6 Lapeyche, Paris, fg St-Honoré, 118.  
 8 Entrée r. Jasmin, 17.  
 10 à 14 Hoine (B<sup>re</sup>), Paris, r. Victoire, 63.  
 16 Audouval et Welry, Paris, r. Raffet, 16.

#### R.8 RAGUINOT (Passage) <sup>18</sup>

- 1 Entrée r. Chalon, 24.  
 3 Lebauf, Alfortville (Seine).  
 5 Robillard, Paris, av. Parmentier, 7.  
 7 Dubois, gér. Rex-Combes, Paris, r. Chaligny, 15.  
 9 Besneux, Paris, r. Lyon, 46.  
 9<sup>bis</sup> Maurel, gér. Salles, Paris, pass. Raguinot, 18.  
 11 Combaco, Paris, av. Grés-Armée, 87.  
 11<sup>bis</sup> Bourdat, Paris, r. Charenton, 100<sup>bis</sup>.  
 13 Buzenet (Vve), Houilles (S.-et-O.).  
 15 Lavohallière, gér. Sagout, Paris, pass. Raguinot, 15.  
 17 Carrier, Pass. Raguinot, 17.  
 19 Lemoine, Chateau Thierry (Aisne).  
 21-23 Jacquelin, Nogent-s-Marne (Seine).  
 25 Leverrier, Paris, pass. Raguinot, 25.  
 27 Maillochon, Paris, r. Levoisier, 183.  
 29 Cheron, gér. Flajol, Paris, pl. Nation, 14.

## Pipeline adapté selon les niveaux de structuration des textes



### R.1 RABELAIS (Rue) <sup>30,31</sup>

- 1 Eiffel, Paris, r. Rabelais, 1.  
 3 Richelot et Ostheimer, Paris, r. Rabelais, 3.  
 5 Entrée r. Montaigne, 26.  
 2-4 Gérard (B<sup>re</sup>), Paris, r. Rabelais, 2.  
 6-8 Aligre de Préaulx (M<sup>re</sup> d'), Paris, fg St-Honoré, 89.

Macrostructure

Microstructure

Annuaire de propriétaires et de propriétés de Paris et de la Seine (Projet Time Machine) 10

## Pipeline adapté selon les niveaux de structuration des textes : Dictionnaire d'Alcedo, entrée pour "Pacajes" (TopUrbi)

del Cerro de Mojanda, y de ella sale un brazo que es el rio Blanco, á la parte del E hay una hacienda que llaman Caxas.

Un rio de la Provincia y Gobierno de Veragua en el Reyno de Tierra-Firme nace en las sierras de Guanico, á la parte del Sur, y desemboca al mar Pacifico.

Otro rio llamado por sobrenombre de los Paeces en la Provincia y Gobierno de Buenos Ayres corre al O, y entra en el de Jacagua, entre el de Joseph Diaz y el paso del Chileno.

Otro de la Provincia y Gobierno del Chocó en el Nuevo Reyno de Granada nace de una laguna, y poco despues se une con el de Quito que nace de otra, y juntos forman el de Atrato.

Una Isla situada en el estrecho de Magallanes, cerca de la Costa del E, delante del Cabo de Monmouth.

Otra Isla pequeña de la mar del Sur en la Bahia de Panamá, delante del golfo de San Miguel.

PAC, Laguna pequeña de la Provincia y Gobierno de Yucatán.

Tiene el mismo nombre un rio pequeño de la Provincia y Gobierno de Guayana ó Nueva Andalucía, nace en el pais de los Indios Caribes feroces, y entra en el Caroni poco despues del parage donde le entra el caudaloso Arui.

PACABARA, Rio de la Provincia y Gobierno de Moxos en el Reyno de Quito, corre al N, y entra en el de Beni.

PACAIPAMPA, Pueblo de la Provincia y Corregimiento de Piura en el Perú, anexo al Curato de Frias.

PACÁJAS, Rio pequeño del pais de las Amazonas, corre al N, entre los de Jacunda y Guanapú, y entra en el Marañon en el brazo que forma la Isla de Joanes; da el nombre este rio á una nacion de Indios que está poco conocida, y habita en la orilla boreal del Marañon, casi 80 leguas mas arriba del Paranaiba.

PACAJES, Provincia y Corregimiento del Perú, confina con la de Chucuito por el NO, por el N con la gran laguna Titicaca, por el NE con la Provincia de Omasuyos, y si-

PACAJES, Provincia y Corregimiento del Perú, confina con la de Chucuito por el NO, por el N con la gran laguna Titicaca, por el NE con la Provincia de Omasuyos, y siguiendo por el E con la Ciudad de la Paz y Provincia de Cicasica, por el SE con el Corregimiento de Oruro y Provincia de Paria, por el S con la de Carangas, por el SO y O con la jurisdicción de la de Arica, mediando la cordillera..

# Pipeline adapté selon les niveaux de structuration des textes

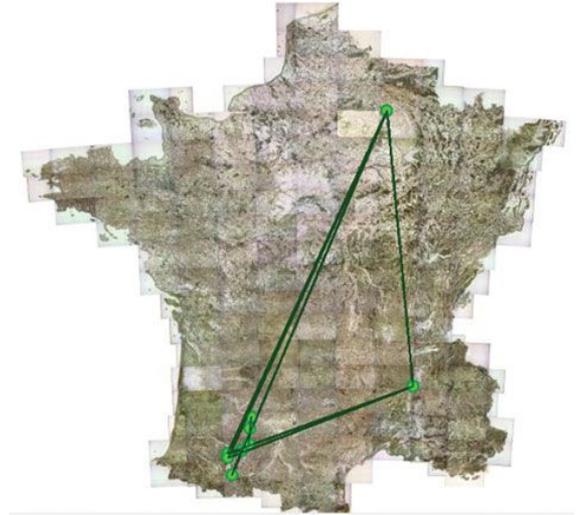
## Extrait sur la Champagne et le Sud-Ouest :

*“Sa puissance émotive, qui fut grande, demeura intacte jusqu'à la fin. Elle partit, le 28, de **Ruffec**, dans les circonstances qu'on connaît. Vers l'Ouest, elle gagna les forêts de **Chizé** et d'**Aulnay**, semble-t-il, à moins que celles-ci n'aient constitué un centre d'émotion locale. ”*

(Paris et al 2017)

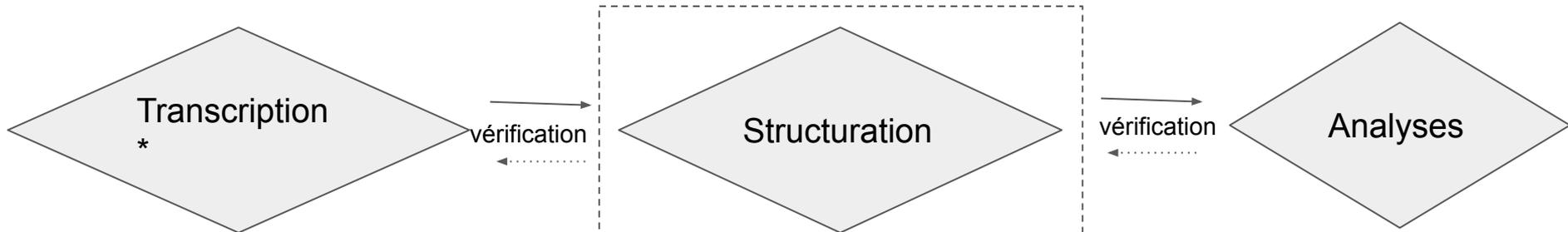
<https://dx.doi.org/10.1080/15420353.2017.1307306>

La dimension spatiale de « La grande peur de 1789 » de Georges Lefebvre

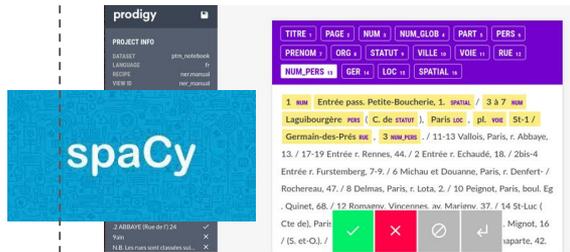


Moins de structure, le cas de noms de lieux est faisable, néanmoins complexe au delà des EN “classiques : relations spatiales, co-références, ...

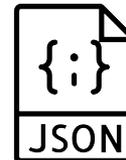
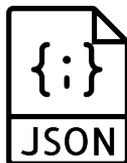
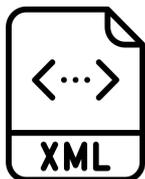
# Une chaîne de traitement et d'analyse : annuaires PTM



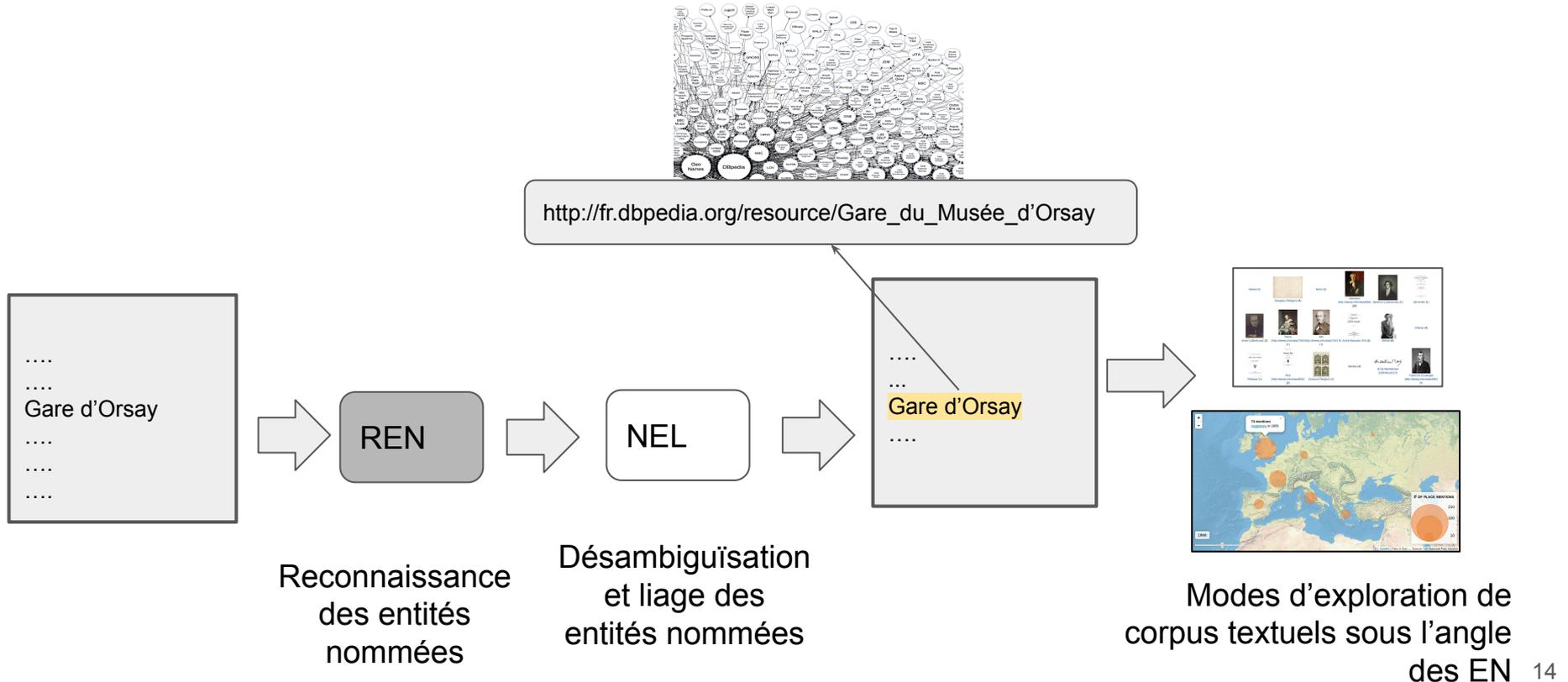
## Outils



## Formats



# Chaîne TAL pour identifier les EN dans des textes littéraires (structuration)



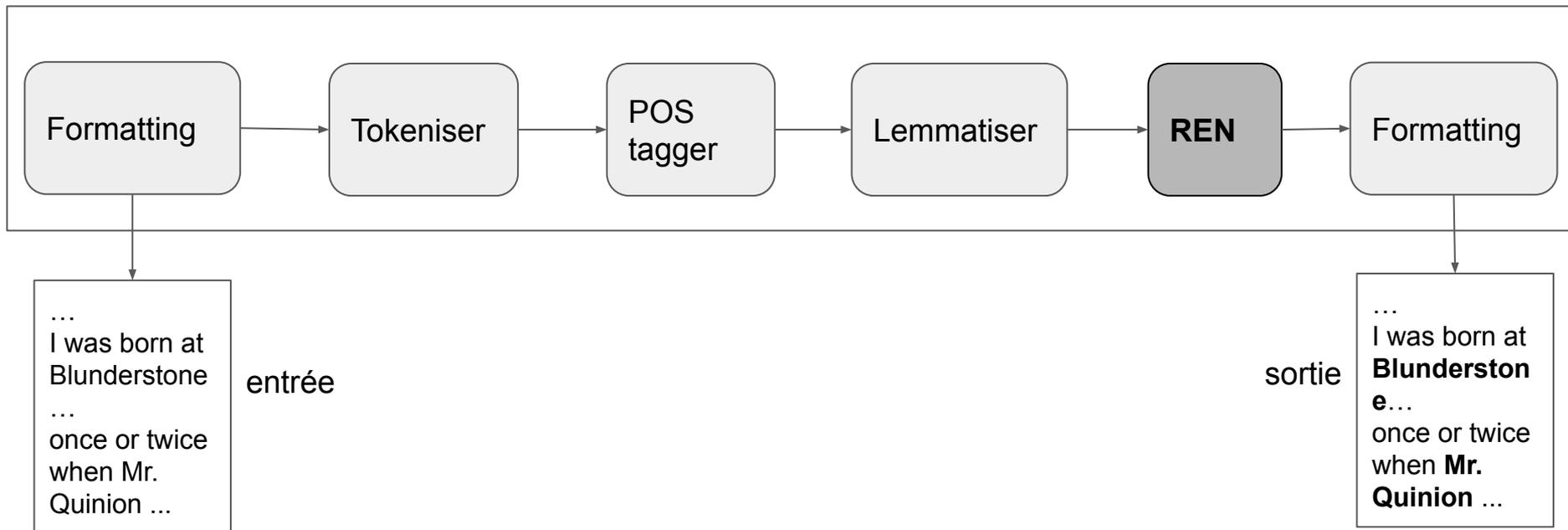
# Enjeux EN et textes anciens : ambiguïtés des localisations

- ❖ imprécision, ("Western Europe")
- ❖ subjectivité ("the main cities of Europe")
- ❖ noms vernaculaires ("Old World" pour l'America)
- ❖ référents changeant de frontières au fil du temps (Allemagne 1850 vs.1950)
- ❖ des lieux historiques qui ont cessé d'exister ("Gaule")
- ❖ ou qui ont changé de nom ("Place de Greve", aujourd'hui place de l'hôtel de ville)

... et plus largement, des référentiels "bruités", peu exhaustifs, des représentations géométriques insuffisantes (ex: pour la Seine)

Exemples issus de : L'Herésiarque et cie (1910) de Guillaume Apollinaire, (23 nouvelles, 54166 mots, fiction) et "Qu'est-ce qu'une nation" (1882) d'Ernest Renan, (8500 mots, essai), corpus OBVIL

# Un pipeline REN typique



Plusieurs configurations sont possibles et, en fonction de la langue, certains composants du pipeline peuvent être optionnels, à l'exception de la tokenisation. La segmentation des phrases peut être un autre composant.

# Annuaire (PTM)

étiquette	description
TITRE	Le titre de la microstructure
NUM	Numéro d'immeuble de la propriété
NUM_GLOB	Numéro (ou ensemble de numéros) référant à un même propriétaire
PART	Part de propriété (constr, terr, usufruit, nue-prop etc) qui est indiqué avant le nom du propriétaire
PERS	Nom du propriétaire
PRENOM	Prénom
ORG	Nom d'une organisation
STATUT	Vient après le nom, peut être Vve, héritiers, conjoints,
VILLE	Ville du domicile.
VOIE	Nom de la voie (rue, rue de, boul etc) de l'adresse personnelle du propriétaire
RUE	Nom de la rue de l'adresse personnelle du propriétaire
NUM_PERS	Numéro d'immeuble de l'adresse perso du propriétaire
GER	pour Gérant et autres, peut contenir Nom et son adresse(rarement)
LOC	Département (si hors Paris) ou pays
SPATIAL	renvoie vers l'entrée sur une autre adresse

## Q.4 QUATRE-VENTS (Rue des) \*\*

- 1 Entrée r. Condé, 2.
- 3 Burnel, Paris, r. Bourdonnais, 36.
- 5 Carreau, Paris, r. St-Denis, 270.
- 7 Tiers, Paris, r. St-Honoré, 91.
- 9 Entrée r. St-Sulpice, 6.
- 11 Entrée r. St-Sulpice, 8.
- 13 Salantin, Paris, r. Quatre-Vents, 13.
- 15 Prévost, Paris, r. Belzunce, 10.
- 17 Morat, Paris, pl. Denfert-Rochereau, 6.
- 19 Lefèvre, Paris, r. Vaugirard, 150.
- 2 Entrée Carrefour Odéon, 14.
- 4 Favre, Valcus-Charnet près Macon (S.-et-L.).
- 6 Guibert, Juvisy (S.-et-O.).
- 8 Cleray, Paris, r. St-Petersbourg, 11<sup>bis</sup>.

## Q.4 QUATRE-VENTS (Rue des) 22

1 Entrée r. Condé, 2 / 3 Burnel, Paris, r. Bourdonnais, 36 / 5 Carreau, Paris, r. St-Denis, 270 / 7 Tiers, Paris, r. St-Honoré, 91 / 9 Entrée r. St-Sulpice, 6 / 11 Entrée r. St-Sulpice, 8 / 13 Salantin, Paris, r. Quatre-Vents, 13 / 15 Prévost, Paris, r. Belzunce, 10 / 17 Morat, Paris, pl. Denfert-Rochereau, 6 / 19 Lefèvre, Paris, r. Vaugirard, 150 /

## Dictionnaire d'Alcedo : catégories à 2 niveaux

	▼ Concept
	▶ Feature Type
	Social Category
	Title-Role
	Error
	▼ GeoEntity
	Landmark
	Settlement
	Territory
	▶ Lemma
	▼ NonGeoEntity
	Organization
	Person
	Human Group
	Qualifier

§§§id\_00442§§§

ALLAUCA, Pueblo de la Provincia y Corregimiento de Yauyos, en el Perú, anexo al

Curato de Tauripampa.

§§§id\_00443§§§

ALLCA, Provincia antigua del Reyno del Perú, al Poniente del Cuzco: estos Indios

bárbaros hicieron una grande y porfiada resistencia á Manco Capac, IV. Emperador de

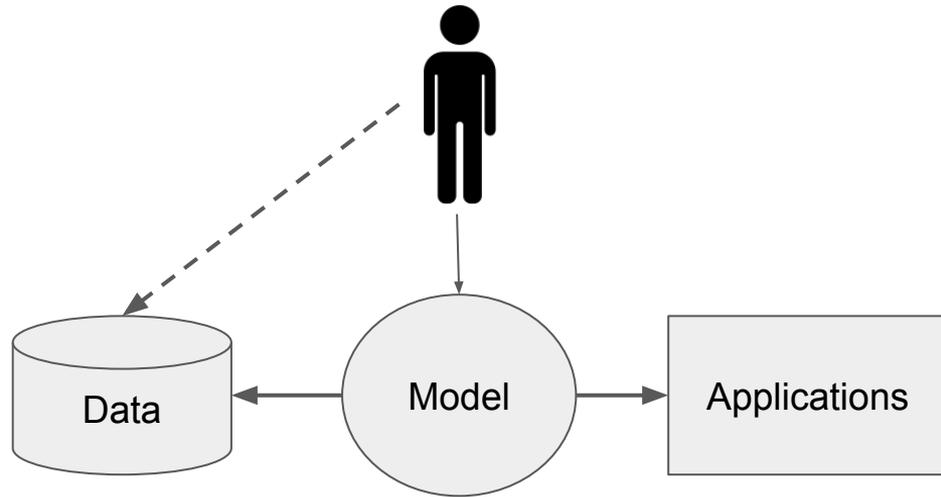
los Incas, llamado el Conquistador, favorecidos de la aspereza del terreno, que

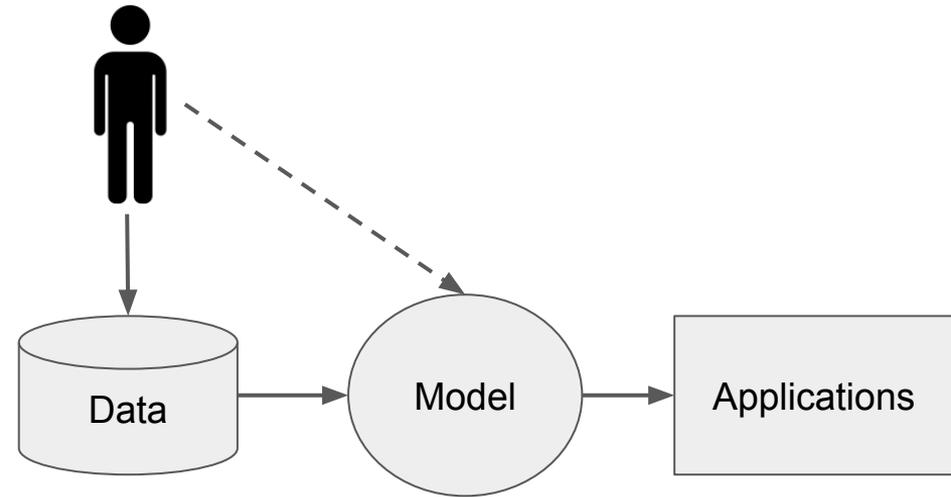
abunda de bosques, montes y lagos, como también de minas de oro y plata.

# De systèmes à base de règles aux approches fondées sur les statistiques et les données



Symbolic systems



Statistical and data-driven approaches

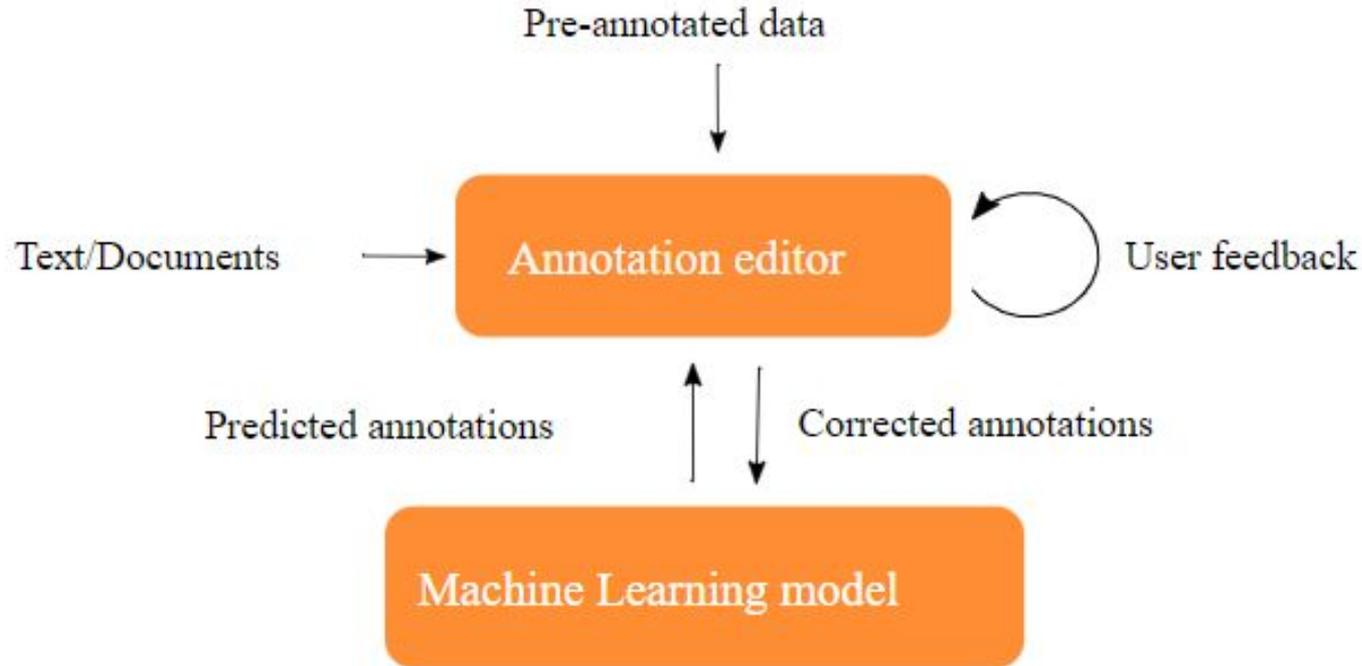


interacts

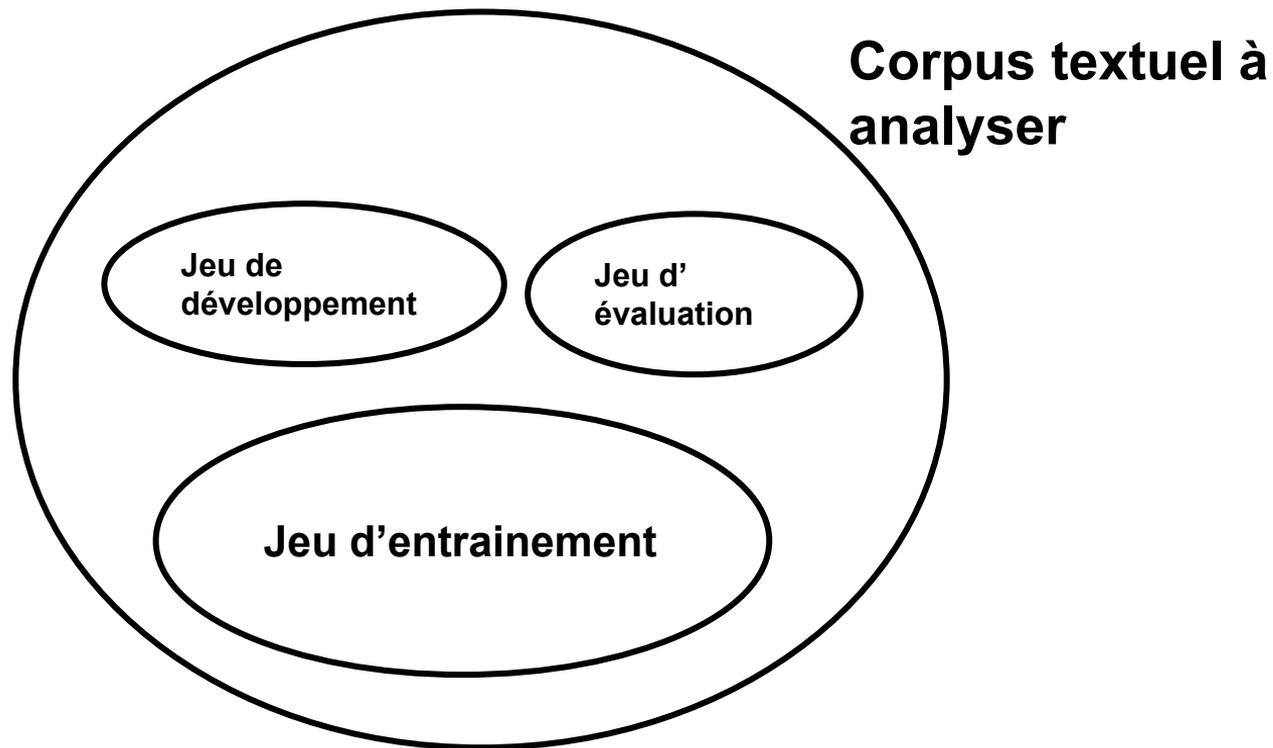


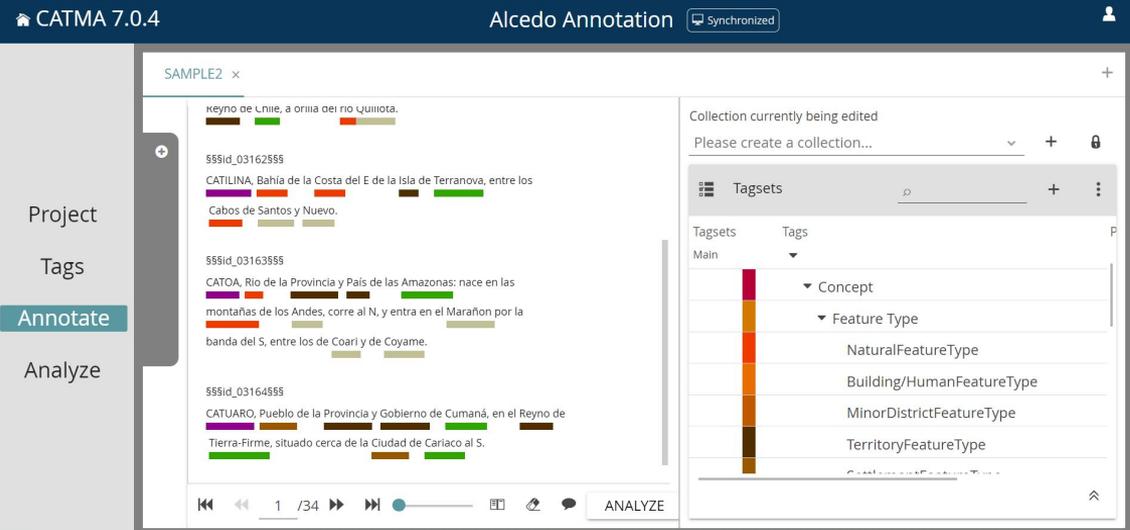
visualise, evaluate, configure

# Approche itérative



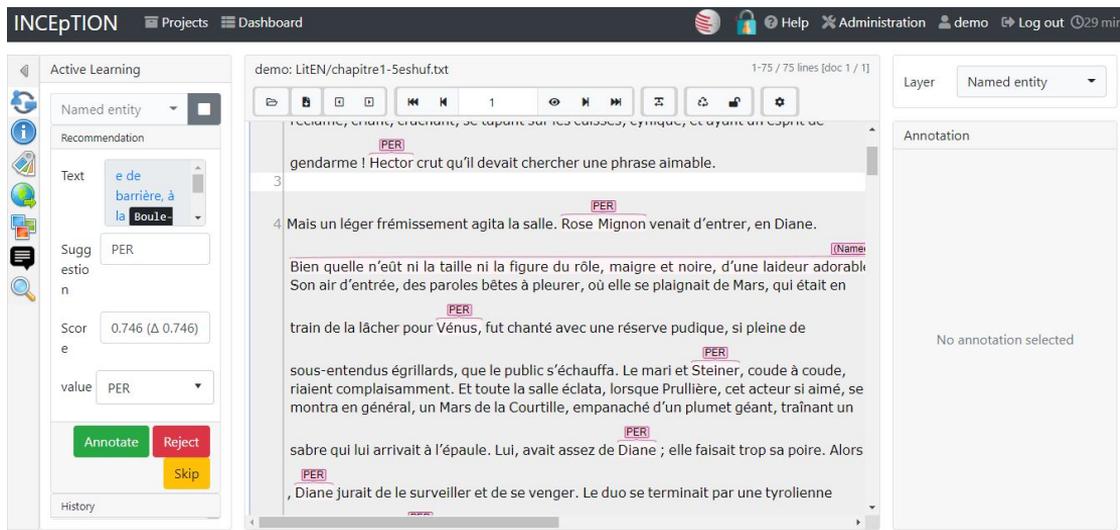
# Adaptation d'un système REN : apprentissage à partir d'un corpus annoté par un utilisateur





# Création du jeu annoté manuellement pour l'entraînement d'un modèle

Catma:  
<https://app.catma.de/seven/>



Inception:  
<https://inception-project.github.io/>

# Evaluation d'un système REN : annuaires PTM

Performance de l'annotation automatique à partir du modèle NER entraîné.

P (précision)  
R (rappel)  
F (f-score)

	P	R	F
NUM	98.68	100.00	99.34
PERS	98.61	99.65	99.12
VILLE	98.96	98.96	98.96
VOIE	99.63	100.00	99.81
RUE	100.00	100.00	100.00
NUM_PERS	100.00	99.62	99.81
LOC	100.00	100.00	100.00
PRENOM	100.00	60.00	75.00
GER	100.00	100.00	100.00
SPATIAL	100.00	98.75	99.37
STATUT	96.67	98.31	97.48
ORG	95.65	88.00	91.67
TITRE	100.00	100.00	100.00
PAGE	100.00	100.00	100.00
NUM_GLOB	100.00	33.33	50.00
PART	100.00	75.00	85.71

Travail réalisé avec Frédérique  
Mélanie-Becquet (LATTICE)

Alrahabi et al, 2021, DOI :  
<https://doi.org/10.4000/revuehn.1079>  
<https://github.com/cvbrandoe/ArabicLitTAL>

[https://office.clarin.eu/v/CE-2020-1738-CLARIN2020\\_ConferenceProceedings.pdf](https://office.clarin.eu/v/CE-2020-1738-CLARIN2020_ConferenceProceedings.pdf)  
<https://github.com/COST-ELTeC>

	Mesure	Partielle	Exacte
<b>Farasa</b>	Précision	<b>0,54</b>	<b>0,42</b>
	Rappel	0,45	0,35
	F1	0,49	0,38
<b>Stanza</b>	Précision	0,51	0,41
	Rappel	0,72	0,58
	F1	<b>0,59</b>	<b>0,48</b>
<b>Madamira</b>	Précision	0,35	0,25
	Rappel	<b>0,55</b>	<b>0,39</b>
	F1	0,43	0,31

<https://github.com/OBVIL/Entites-nommees>  
 (préliminaires, Zola, evaluation L3I à venir)

	F1 (exacte)	F1 (partielle)
Stanza	<b>0,69</b>	<b>0,55</b>
Spacy	0,65	0,40

	Cat	Precision	Recall
SEM-fra	LOC	0.465	0.395
	PERS	0.416	0.138
SPACY-fra	LOC	0.180	0.569
	PERS	0.526	0.629
PALAVRAS-por1	LOC	0.835	0.780
	PERS	0.905	0.901
SPACY-por1	LOC	0.338	0.728
	PERS	0.554	0.645
PALAVRAS-por2	LOC	0.624	0.693
	PERS	0.750	0.866
SPACY-por2	LOC	0.284	0.734
	PERS	0.591	0.707
Stanford-eng	LOC	0.438	0.495
	PERS	0.619	0.548
SPACY-eng	LOC	0.366	0.495
	PERS	0.691	0.453
SrpNER-srp	LOC	0.849	0.578
	PERS	0.820	0.729
SPACY-srp	LOC	0.354	0.308
	PERS	0.637	0.561

Évaluer des systèmes REN prêts à utiliser pour mesurer leur faiblesses avec l'idée de créer un benchmark.

# Approches statistiques et fondées sur les données

La manière dont les données sont soumises au système est cruciale car :

- la **quantité** et la **qualité** des données peuvent rendre le système plus ou moins précis (peu de fausses correspondances), complet (peu de correspondances manquées), robuste (résistance au bruit),
- Le **type** de texte sur lequel l'entraînement est effectué conditionne l'applicabilité du modèle à d'autres types de texte,
- Le **prétraitement** (tokenisation, ...) et la façon dont les entités nommées sont définies influencent la façon dont elles sont reconnues.



# Lier les EN aux référentiels

**Named Entity Linking**, il s'agit d'un moyen d'établir un lien explicite entre les mentions d'EN cités dans le texte et les objets du monde auxquels elles réfèrent. Créer ce lien est important afin de lever toute ambiguïté sur l'identité de l'EN. Il y a donc deux objectifs :

- ❑ Désambiguiser

"Goncourt" - Edmond de Goncourt ou Jules de Goncourt ?

"Voltaire", "François-Marie Arouet" - manières de désigner la même personne

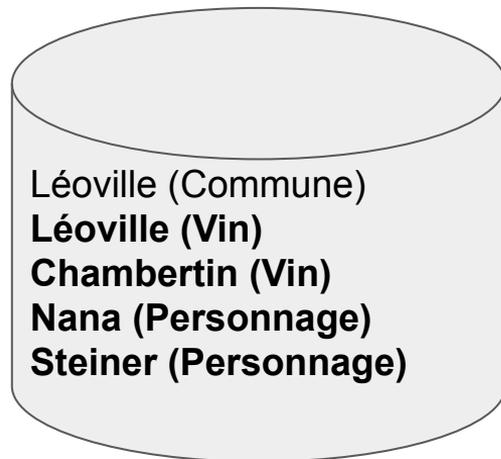
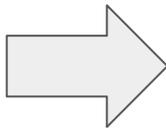
- ❑ Lier à un référentiel

"Voltaire" - associer l'entrée correspondante dans Wikidata  
(<https://www.wikidata.org/wiki/Q9068>)

# Lier les EN aux référentiels

**Léoville** ou **Chambertin** ?  
murmura un garçon, en allongeant  
la tête entre **Nana** et **Steiner**, au  
moment où celui-ci parlait bas à  
la jeune femme.

Texte annoté en EN :  
Nana (Zola)



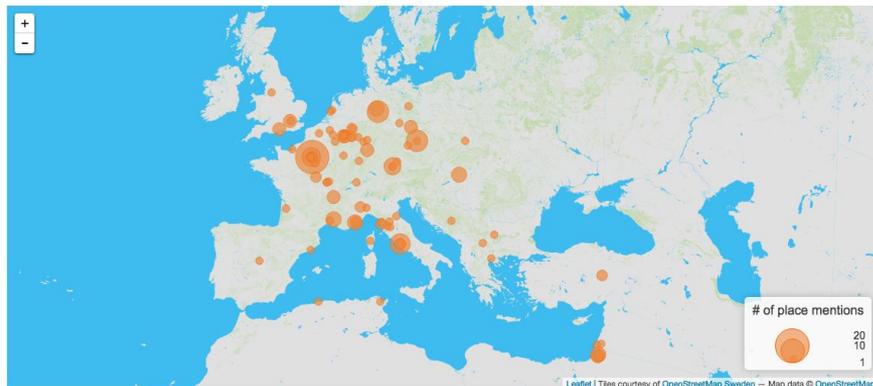
Référents  
possibles

La mention étudiée hors de son contexte est ambiguë. Un humain peut typiquement lever l'ambiguïté par utilisation conjointe de **la connaissance du monde** dont il dispose et **d'indices contextuels**.

# L'empreinte spatiale d'un texte pour rendre compte de la géographie d'un auteur et les échelles mobilisées

## Visualization

The following map shows distribution of places mentioned in the input TEI-XML file, geo-coordinates are obtained via French DBpedia.



140 places are displayed on the map.

19 places were not included on the map because geo-coordinates were unavailable, these are: Rhin, Hambourg, Queensland, Bohême, La Nouvelle-Orléans, Berlin, Amsterdam, Neckar, Ile-de-France, Bavière, montagnes Rocheuses, Provence, Monte-Carlo, empire des Habsbourg, Moldau, Danube, Europe, royaume de Juda, Savoie

## Visualization

The following map shows distribution of places mentioned in the input TEI-XML file, geo-coordinates are obtained via French DBpedia.



137 places are displayed on the map.

17 places were not included on the map because geo-coordinates were unavailable, these are: Berlin, Rhin, Bohême, montagnes Rocheuses, Queensland, royaume de juda, La Nouvelle-Orléans, Provence, Neckar, empire des Habsbourg, Danube, Moldau, Ile-de-France, Hambourg, Bavière, Savoie, Amsterdam

You can download the resulting annotated XML-TEI file [here](#)

Guillaume Apollinaire's « Le passant de Prague »

# Quelques remarques finales : enjeux EN et corpus des HN

- **langues de textes** diverses, anciennes, mélangées, non normalisés
- **catégories EN** fines vs. trop générales ; différentes **définitions des EN** selon les cas d'applications (ex: édition TEI)
- Intégrer dans les modèles REN :
  - des éléments de la **mise en forme** du document imprimé (travaux de Joseph Chazalon et l'EPITA et l'IGN sur les annuaires dits Bottin)
  - **bruit** présente dans le texte issu de la numérisation (voir Séminaires du Groupe NER for historical documents <https://ner-for-historical-docs.github.io/>)
- **référentiels** de données anciennes peu disponibles pour lier les EN
- **évaluation** des performances des systèmes : manque de jeux de tests pour comparer
- **Modèles à adapter**, fine-tuning, modèles génératifs, besoin de ressources de calcul
- **Au delà des EN**: annoter les vocabulaires contrôlés et expressions polylexicales reste simple, les relations spatiales faisable dans la plupart des cas, mais peu de propositions pour les entités non nommées et les co-références présentes dans des textes littéraires (sauf pour les travaux en cours du laboratoire LATTICE sur le French BookNLP)

# Ressources pour découvrir le REN

Spacy - modèles disponibles : <https://spacy.io/models>

Quelques notebooks Python pour démarrer :

1) <https://github.com/cvbrandoe/coursTAL/tree/master/2023/notebooks> :

coursTAL\_EltecFRL2\_StepA\_1eAnnotStanza\_CorrTagTog.ipynb

coursTAL\_EltecFRL2\_StepB\_TrainSpacy.ipynb

coursTAL\_EltecFRL2\_StepC\_UseTrModelSpacy.ipynb

2) <https://github.com/ludovicmoncla/tutoriel-geoparsing>

3) <https://github.com/orgs/Consortium-ARIANE/>

(issus de la formation ARIANE à Lyon, ils seront déposés ici la mi-novembre 2023)

MERCI

# Evaluation d'un système REN

gold: *Phébus parut en Postillon de Lonjumeau et Minerve en Nourrice normande.*

test: *Phébus parut en Postillon de Lonjumeau et Minerve en Nourrice normande.*

True positive (TP)   False Negative (FN)   TP   False positive (FP)

- ❖ **Recall** : the number of correctly annotated entities wrt the total of manually annotated entities in the gold  
$$= TP / (TP + FN)$$
- ❖ **Precision** : the number of correctly annotated entities wrt the total of returned entities  
$$= TP / (TP + FP)$$
- ❖ **F-score**, harmonic mean between recall and precision.

It is also possible to **extend** these measures to use **relaxed match** when comparing annotations, instead of **strict match** as presented above.